

Europäisches Patentamt
European Patent Office
Office européen des brevets



(11) **EP 0 844 561 A2**

(12) **EUROPEAN PATENT APPLICATION**

(43) Date of publication:
27.05.1998 Bulletin 1998/22

(51) Int. Cl.⁶: **G06F 11/10, G11B 20/18**

(21) Application number: **97120114.0**

(22) Date of filing: **17.11.1997**

(84) Designated Contracting States:
**AT BE CH DE DK ES FI FR GB GR IE IT LI LU MC
NL PT SE**
Designated Extension States:
AL LT LV MK RO SI

- Ogata, Mikito
Odawara-shi (JP)
- Kurano, Akira
Odawara-shi (JP)
- Tamiya, Toshihiko
Hadano-shi (JP)
- Yamamoto, Akira
Sagamihara-shi (JP)
- Takahashi, Naoya
Yokohama-shi (JP)

(30) Priority: **21.11.1996 JP 310520/96**

(71) Applicant: **Hitachi, Ltd.**
Chiyoda-ku, Tokyo 101-0062 (JP)

(72) Inventors:
• Katsuragi, Eiжу
Odawara-shi (JP)

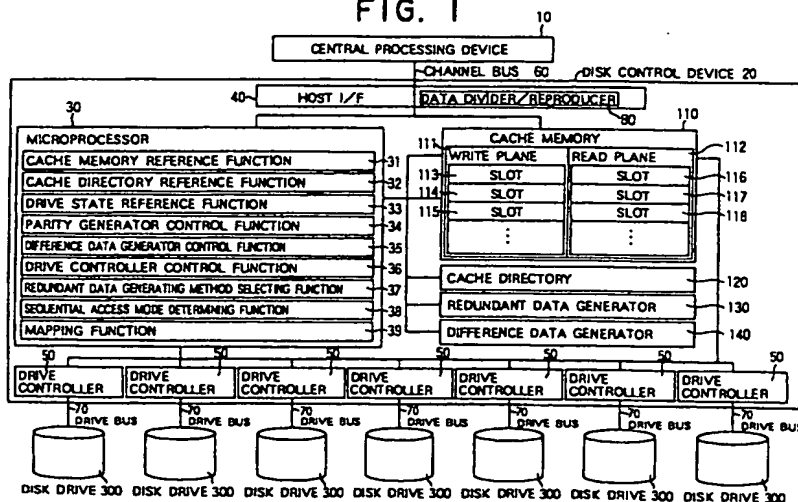
(74) Representative:
Strehl Schübel-Hopf & Partner
Maximilianstrasse 54
80538 München (DE)

(54) **Disk array device and method for controlling the same**

(57) The disk array device includes a disk control device (20) connected to a central processing unit (10) and a plurality of disk drives (300) composing disk arrays under the control of said disk control device (20). The disk control device (20) includes a redundant data generator (130), a difference data generator (140), and a redundant data generation method selecting function (37). The disk array device selects a proper redundant data generating method from a method of read and modify and a method of all stripes, both of which are executed to generate redundant data by the disk control

device (20) according to an access pattern from a host, a load state of the disk drive (300), and a failure, and a method of a generation in a drive and a method of difference, both of which are executed to generate the redundant data on the disk drive (300) for saving the redundant data, for the purpose of reducing an overhead accompanied with generation of the redundant data and improving reliability of generating the redundant data.

FIG. 1



EP 0 844 561 A2

Description

BACKGROUND OF THE INVENTION

5 Field of the Invention

The present invention relates to a disk array technique and a technique for controlling the disk array, and more particularly to the technique for effectively enhancing efficiency and reliability of a process for redundant data generated in writing data.

10

Description of the Related Art

David A. Paterson, et al. have reported in UCB/CSD/87.391 (December 1987) of University of California's Report, a method for saving redundant data in part of a storage device. This method takes the steps of preparing a plurality of storage devices (disk drives), dividing the I/O data from a host computer into the corresponding number of parts to the number of the storage devices, recording and reproducing the divided data into and from the storage devices, and recovering fault data from a normal storage device(s) if a temporary or permanent failure takes place in some of the storage devices. This reports said that the following two methods are provided for generating redundant data.

The first method for generating the redundant data, which is referred to as a method of read and modify, is arranged to use write data from a host computer, previous data (data before update) stored in a storage device where the write data is to be stored, and previous data (data before update) stored in a storage device where the generated redundant data is to be stored for the purpose of generating the redundant data. In this method, assuming that the divisional number of the write data is a , for all the storage devices, the I/Os for writing take place $a + 1$ times and the I/Os for reading take place $a + 1$ times as well. The total I/O times reach $2 \times a + 2$. If the storage device does not contain the redundant data, the access times are a . It means that the use of the redundant data results in increasing the I/O times by $a + 2$.

The second method for generating the redundant data, which is referred to as a method of all stripes, is arranged to use the divided parts of write data from the host computer and data read from the storage devices except those for saving the redundant data belonging to the ECC group that does not save the write data from the host computer, for the purpose of generating the redundant data.

With this method, assuming that the divisional times of the write data is a and the number of the storage devices composing an ECC group, except those for saving the redundant data, is b , the times of I/Os to and from the storage devices are totally $b + 1$, wherein the I/O times for reading are a and the I/O times for writing data containing the redundant data are $a + 1$. For the storage devices that do not contain the redundant data, since the access times are a , the use of the redundant data results in increasing the I/O times by $b + 1 - a$.

Apart from the foregoing methods, the method for generating the redundant data in the storage device has been disclosed in a U.S. Patent No. 5,613,088, wherein the redundant data is generated in the storage device provided with two heads for read and write. Concretely, the read head and the write head are fixed on a common actuator so that the read head reads parity data before update and then the write head writes on the same area updated parity data generated from the parity data before update unless the disk spins once.

The foregoing two methods for generating the redundant data, that is, the method of read and modify and the method of all stripes, have the increased number of I/Os when writing the data because of the use of the redundant data. This means that the disk control device with redundancy is inferior in performance to the disk control device without redundancy. Hence, the conventional disk control device with redundancy selectively employs the method with a smaller number of I/Os to and from the storage device for the purpose of reducing the I/Os to and from the storage device in writing data. This selection makes it possible to reduce the burden on the storage device and thereby improve the processing speed. Concretely, in the case of $a \geq (b - 1)/2$, the method of all stripes has a smaller increase of I/O times for the storage device than the method of read and modify, while in the case of $a < (b - 1)/2$, the method of read and modify has a smaller increase. By using this, if the data length of the write data from the host computer is in the range of $a < (b - 1)/2$, for example, in the case of a transaction process, the disk control device arranged to use the method for generating the redundant data operates to generate a parity through the effect of the method of read and modify. The I/O times for the storage device becomes four at minimum when $a = 1$, which is the smallest load burdened on the storage device. In other words, however, in such a case of the transaction process, it means that the critical point of performance takes place when $a = 1$. The performance cannot be improved further unless the method of process at this time is reconsidered. The problem about the method of read and modify is essentially based on the fact that two I/Os are issued to the storage device for saving the redundant data and a mechanical overhead such as movement of the head and spinning on standby is burdened at each I/O time. The mechanical overhead is a great bottleneck on the disk control device that is in electric operation.

EP 0 844 561 A2

The method disclosed in a U.S. Patent No. 5,613,088 makes it possible to generate the redundant data in the storage device provided with two heads for read and write, thereby reducing the spinning times of the drive on standby. In case this method is expanded to a general storage device provided with a single head, the resulting method takes the steps of transferring the updated data and the data before update read from the storage device to the storage for saving the redundant data and enabling the storage device for saving the redundant data to read the redundant data before update and generate the redundant data from the transferred updated data, the data before update, and the redundant data before update. This method is referred to as a method of a generation in a drive. In this method, when the head is positioned to read the redundant data before update after the spinning on standby and reaches the next writing position, the write is started. This operation makes it possible to avoid the spinning on standby during the writing interval and merely needs one movement of the head and one standby spin. As a result, if the length of the data from the host computer is short, the processing speed of the control device can be improved further.

However, the method of a generation in a drive cannot necessarily offer the essential effects depending on an access pattern from the host computer, the load burdened on the storage device, and the like.

That is, if the length of the generated redundant data is longer than one spin of the disk, that is, the length of the write data from the host computer is longer than one spin of the disk, the method of a generation in a drive is required for the disk to spin on standby during the interval of reading the redundant data before update and writing the updated redundant data. However, the spinning on standby results in increasing an occupying time of the drive for saving the redundant data, thereby increasing a response time of the drive for saving the redundant data. If, therefore, the length of the redundant data is one spin, the method of a generation in a drive has a problem that the disk array device has a lower response time.

The method of a generation in a drive is arranged to increase the data before update to be read out as the divisional number of the write data from the host computer becomes more, thereby increasing the load burdened on the storage device. Hence, if the divisional number of the write data is great, the method of a generation in a drive disadvantageously puts the throughput of the disk array device into a lower value.

The use of the method of a generation in a drive makes it possible to increase an occupying time of the drive for saving the redundant data at each spin as compared with the method of read and modify, thereby increasing the load burdened on the drive for saving the redundant data in the highly multiplexing and high load environment. Hence, the method of a generation in a drive may enhance a probability that the drive for saving the redundant data is in use, thereby lowering the throughput of the drive.

When the write data is transferred from the host computer to the disk control device together with an explicit specification of consecutive pieces of data, the method of a generation in a drive operates to immediately generate the redundant data on the transferred write data. As a result, when the succeeding write data is transferred from the host computer, the method of all stripes may lose a chance of generating the redundant data in correspondence to the first write data. Hence, if the method of a generation in a drive cannot use the method of all stripes, this disadvantageously lowers the efficiency of generating the redundant data, thereby degrading the throughput of the disk array device.

When the method of a generation in a drive uses the generation of the redundant data, the generation of the redundant data becomes unsuccessful because of any failure such as failure caused in reading the redundant data before update. In this case, the redundancy of the ECC group may be lost at once.

SUMMARY OF THE INVENTION

It is an object of the present invention to avoid the spinning on standby caused if the data length of the write data from an upper system is longer than one spin and improve a response time of the disk array device caused if the disk drive composing the disk array generates the redundant data.

It is a further object of the present invention to improve a throughput of the disk array device by selecting the most approximate method for generating redundant data so that the necessary reading number of the data before update is made minimal according to the divisional number of the write data received by the upper system.

It is a yet further object of the present invention to improve a response time of the disk array device by reducing an occupying time for one spin of the disk drive in the case of generating the redundant data through the effect of the disk drive composing the disk array.

It is another object of the present invention to improve a throughput of a disk array device by enhancing the efficiency of generating the redundant data in association with the process for write data according to an access pattern required by the upper system.

It is still another object of the present invention to enhance reliability of a process for generating the redundant data through the effect of the disk drive composing the disk array.

According to the invention, a disk array device having a plurality of disk drives composing a disk array and a disk control device for controlling those disk drives includes a plurality of means for generating the redundant data in different ways, and a selective control logic for selectively executing at least one of the plurality of means for generating

EP 0 844 561 A2

redundant data.

The disk array device including a plurality of disk drives composing a disk array and a disk control device for controlling the disk drives is dynamically switched from the generation of the redundant data in the disk control device to the generation of the redundant data inside of the disk drive according to an operating status.

Concretely, as an example, the disk array device according to the invention includes the following components.

That is, the disk array device composing a plurality of disk drivers, which is arranged to set a logic group of a partial set of the disk drives and save the redundant data in part of the logic group for the purpose of recovering the fault data from the normal disk drives when some of the disk drives are disabled by temporary or permanent failure, provides means for generating the redundant data in each of the disk drives.

The disk control device includes a first redundant data generating circuit inside of a control device for generating new redundant data from the divided write data received from the upper system, and the previous data to be updated by the divided write data, and the redundant data of the data group of the divided write data, a second redundant data generating circuit inside of the control device for generating new redundant data of the data group from the data that is not updated by the divided write data contained in the data group, and a selective control circuit for selecting a difference data generating circuit for generating difference data from the divided write data received from the upper system and the previous data updated by the divided write data and the circuit for generating the redundant data.

Further, the disk control device provides means for determining a data length of the write data received from the upper system, means for detecting a load of the drive for saving the redundant data, means for determining if the transfer of consecutive pieces of data from the upper system is explicitly specified, and means for detecting if the generation of the redundant data inside of the disk drive for saving the redundant data is failed. The selective control circuit operates to select a proper method for generating the redundant data.

The disk array device and the method for controlling the disk array device as described above are served as follows an example.

The disk array device operates to determine the data length of the write data sent from the upper system and generate the redundant data in the disk drive if the data length is determined to be shorter than one spin of the disk. Hence, the disk array device operates to suppress the spinning on standby caused inside of the disk drive for saving the redundant data, thereby improving a throughput of the disk array device.

If the data length of the write data sent from the upper system is determined to be longer than one spin of the disk, the difference data between the divided write data and the data before update stored on the disk drive for saving the divided data is transferred onto the disk drive for saving the redundant data. The disk drive for saving the redundant data operates to generate the redundant data from the difference data and the redundant data before update, thereby suppressing the spinning on standby caused in the disk drive for saving the redundant data and improving the throughput of the disk array device accordingly.

The method for controlling the disk array device is executed to determine a load burdened on the disk drive for saving the redundant data, generate the redundant data in another disk control device without having to actuate the method of a generation in a drive if the load is determined to be greater than or equal to a given value, for the purpose of distributing the load associated with the generation of the redundant data. In the highly multiplexing and high load environment, by suppressing the increase of the load to be put on the disk drive for saving the redundant data, it is possible to suppress the probability that the disk drive for saving the redundant data is in use and thereby improve the throughput of the disk array device.

The method for controlling the disk array device is executed to determine if the transfer of consecutive pieces of data from the upper system to the disk control device is explicitly specified and generate the redundant data in the block a short time after the write data reaches a sufficient length without immediately generating the redundant data, if the explicit transfer of the consecutive data is specified. This enables to improve the efficiency of generating the redundant data and the throughput of the disk array device.

When the method of a generation in a drive fails in generating the redundant data, the method for generating the redundant data is switched to another method. This makes it possible to increase the changes of recovering the failure and thereby improving reliability of the disk array device.

BRIEF DESCRIPTION OF THE DRAWINGS

Fig. 1 is a diagram showing an arrangement of an information processing system including a disk array device according to the present invention;

Fig. 2 is a view showing an internal arrangement of a disk drive used in the disk array device according to the present invention;

Fig. 3 is a view showing an example of mapping data to be given to and received from an upper system in the disk drive used in the disk array device according to the present invention;

Fig. 4 is a diagram showing an arrangement of hardware of an information processing system including a disk array

EP 0 844 561 A2

device according to the present invention;

Fig. 5 is a flowchart showing an example of a process for selecting a method for generating redundant data in a disk array device according to the present invention;

Fig. 6 is a concept view showing a flow of data executed in generating redundant data through the effect of a method of all stripes in the disk array device according to the present invention;

Fig. 7 is a concept view showing a flow of data executed in generating redundant data through the effect of a method of a generation in a drive in the disk array device according to the present invention;

Fig. 8 is a concept view showing a flow of data executed in generating redundant data through the effect of a method of read and modify in the disk array device according to the present invention;

Fig. 9 is a concept view showing a flow of data executed in generating redundant data through the effect of a method of difference in a disk array device according to an embodiment of the present invention; and

Fig. 10 is a view showing an arrangement of a cache directory in the disk control device included in the disk array device according to an embodiment of the present invention.

15 DESCRIPTION OF THE PREFERRED EMBODIMENTS

Hereafter, an embodiment of the invention will be described in detail with reference to the appended drawings.

Fig. 1 is a concept showing an arrangement of an information processing system including a disk array device according to an embodiment of the present invention. The information processing system according to this embodiment

is arranged to have a central processing device 10 (referred to as a host 10) and a disk array device connected thereto. The disk array device is configured of a disk control device 20 and seven disk drives 300 (300a to 300g) to be operated independently of each other, each of which disk drive may be a magnetic disk unit, for example. These seven disk drives 300 compose an ECC group (on which unit data is recovered when failure takes place. The host 10 is coupled with the disk control device 20 through a channel bus 60. The disk control device 20 is connected to each disk drive 300 through

the corresponding drive bus 70 so that each disk drive 300 may be operated independently of each other. The disk control device 20 is arranged to have a host I/F 40, a data divider/reproducer 80, a microprocessor 30, a cache memory 110, a cache directory 120, a redundant data generator 130, a difference data generator 140, and a drive controller 50. The host I/F 40 and the data divider/reproducer 80 is coupled to a channel bus 60, the cache memory 110, and the microprocessor 30 through signal lines. The microprocessor 30 is coupled to the cache memory 110, the cache directory 120, the redundant data generator 130, the difference data generator 140, and the drive controller 50 through signal lines. The cache memory 110, the cache directory 120, the redundant data generator 130, the difference data generator 140, and the drive controller 50 are controlled by a cache memory reference function 31, a cache directory reference function 32, a drive state reference function 33, a redundant data generator control function 34, a difference data generator control function 35, and a drive controller control function 36, all of which are executed by

micro programs built in the microprocessor 30. The microprocessor 30 includes a redundant data generating method selecting function 37 served as means for determining the method for generating the redundant data, a sequential access mode determining function 38 served as means for checking if the sequential access from the host 10 is specified, and a mapping function 39 served for calculating a write position on the actual disk drive 300 from the write data from the host 10. Those functions are executed

by the micro programs built in the microprocessor 30 itself. The cache memory 110, the cache directory 120, the redundant data generator 130, and the difference data generator 140 are coupled through signal lines. The cache memory 110 is coupled to the drive controller 50 through a signal line so that data is allowed to be transferred between the cache memory 110 and the drive controller 50. The cache memory 110 is divided into a write plane 111 on which the write data from the host 10 is saved and a read plane 112 on which the data read from the disk drive 300 is saved. Each of the read plane 112 and the write plane 111 contains slots 113 to 118 on which each plane is divided into sectors.

Fig. 10 is a concept view showing an arrangement of a cache directory 120 according to this embodiment. In this embodiment, the cache directory 120 contains several kinds of information to be set thereto. Those kinds of information include cache managing information 121 for managing data of the cache memory 110, data length information 122 for saving a data length received from the host 10, access pattern information 123 for storing the access pattern in the case of specifying an access pattern such as a random access and a sequential access from the host 10, and pending flag information 124 for saving a pending flag for holding a data write process containing generation of redundant data until a series of sequential accesses are terminated if the access pattern is the sequential pattern.

Fig. 2 is a concept view showing an internal arrangement of the disk drive used in this embodiment. The disk drive 300 includes a disk I/F 310 for controlling transfer of information between the disk drive 300 and the outside through the drive bus 70, a drive buffer 330 for temporarily holding data received from the outside and data read from an inside disk medium 350, a disk control mechanism 340 for controlling a positioning operation of a head (not shown) with respect to the disk medium 350, and a microprocessor 320 for controlling all of those components. The disk I/F 310, the drive

EP 0 844 561 A2

buffer 330 and the disk control mechanism 340 are coupled to the microprocessor 320 through signal lines and are controlled by the microprocessor 320.

In this embodiment, the microprocessor 320 of each disk drive 300 has a function of generating new redundant data from the data received from the outside and the redundant data before update saved in the disk medium 350 inside of the disk drive 300 itself and then saving the new redundant data in the disk medium 350 inside of the disk drive 300. This function may be executed by the microprocessor 320 or any leased hardware except the microprocessor 320.

Fig. 3 is a concept view showing an example of mapping I/O data to be given to or received from the host 10 onto the disk medium 350 located inside of the disk drive. In this embodiment, the data recording area of the disk medium 350 of each disk drive 300 is logically divided into a plurality of unit areas. The concatenation of each unit area of the disk drives 300a to 300g composes a data group containing at least one piece of redundant data.

That is, the redundant data P000 (parity) is located on the unit area of the most right disk drive 300g of the first column. On the second column or later, the redundant data is shifted by one to the left hand of the storage position of the parity on the previous column. If the storage position of the redundant data of one previous column is located on the leftmost disk drive 300a, the redundant data P001~ is located on the most right disk drive 300g. The divisions D000 to D029 of the write data from the host 10 are sequentially mapped from the disk drive located immediately in the right hand of the redundant data or the leftmost disk drive 300 if the redundant data is located in the most right hand. The redundant data of each column is generated to be equal to the exclusive OR of the data of each column, for example, D000 to D005 and then is saved. If one failure takes place in one data piece of each row, the fault data can be recovered by the exclusive OR of the remaining data inside of the column with the redundant data.

In addition, the redundant data for error recovery may be an exclusive OR of a group of plural divided data as well as any code such as a hamming code.

Fig. 4 is a concept view showing the arrangement of the disk array device according to this embodiment. The channel bus 60 shown in Fig. 1 corresponds to 260-1 to 260-8. The host I/F 40 and the data divider/reproducer 80 shown in Fig. 1 correspond to host adapters 231-1 and 231-2. The microprocessor 30, the drive controller 50, the redundant data generator 130, and the difference data generator 140 shown in Fig. 1 correspond to disk adapters 233-1 to 233-4. The cache memory 110 shown in Fig. 1 corresponds to cache memories 232-1 to 232-2. The cache directory 120 shown in Fig. 1 corresponds to shared memories 234-1 to 234-2. The drive path 70 shown in Fig. 1 corresponds to 270-1 to 270-16.

The host adapters 231-1 to 232-2, the cache memories 232-1 to 232-2, the disk adapters 233-1 to 233-4, and the shared memories 234-1 to 234-2 are connected to each other through doubled data transfer buses 237-1 to 237-2.

Under the control of the disk control device 20, the inside of each of the disk drive boxes 241-1 to 241-2 is connected to a storage device 240 arranged to accommodate the disk drives 242-1 to 242-32 and the disk drives 242-33 to 242-64.

The host adapters 231-1 to 232-2, the disk adapters 233-1 to 233-4, and the shared memories 234-1 to 234-2 are connected to a service processor 235 through a communication path 236 inside of a control unit. This service processor 235 is handled from the outside through a maintenance terminal 250.

Fig. 1 illustrates a single ECC group (group of disk drives), while Fig. 4 illustrates totally eight ECC groups (group of disk drives). The disk drives 242-1 to 242-7, 242-9 to 242-15, 242-17 to 241-23, and 242-25 to 242-31 correspond to ECC groups. The disk drives 242-8, 242-16, 242-24 and 242-32 are spared disk drives. So are the disk drives 242-33 to 242-64.

The description will be oriented to how the microprocessor is operated with respect to the disk control device 20 arranged as described above when the host 10 issues the write data to the microprocessor along the flow shown in Fig. 5. The write data from the host 10 is transferred to the disk control device 20 through the channel bus 60 and is divided into sector lengths by the data divider/reproducer 80. When the host I/F 40 serves to save each division in the slot of the write plane 111 of the cache memory, at a step 1000 of Fig. 5, the microprocessor 30 operates to count the number of data lengths transferred from the host I/F 40 to the cache memory 110 and then save the count value in a data length information 122 included in the cache directory 120. Proceeding to a step 1010, the sequential access mode determining function 38 is executed to determine if the access specification is sequentially given from the host 10. The presence or absence of the sequential access specification is saved in the host access pattern information 123 of the cache directory 120.

If the sequential access is specified, the operation goes to a step 1150. At this step, the termination of the writing process is reported to the host 10 without immediately generating the redundant data. Then, for indicating the write data is not reflected on the disk drive 300, a pending flag is raised to the pending flag information 124 on the cache directory 120 and then the microprocessor 30 is waiting for a specified time. This is because there exists a high probability that the succeeding write data from the host 10 is issued if the sequential access is specified and that the method of all stripes may be used for generating the redundant data by using plural pieces of write data. Hence, the transfer of the succeeding data allows the microprocessor to just wait until the stripes (a set of slots to be written on the same position of the different disk drives) are made to be the write data. The proper waiting time is calculated by

EP 0 844 561 A2

(date length of the stripe - transferred data length)/data transfer speed from the host. After the wait for a specified time, the operation goes to a step 1020.

If no sequential access is specified from the host at the step 1010, the operation at the step 1020 is immediately executed. The microprocessor 30 operates to calculate the positions of the disk drive 300 where the write data of each slot is saved and of the disk medium 350 through the effect of the mapping function 39 for the purpose of deriving the column on the disk drive 300 where the write data is to be written. Next, the microprocessor 30 operates to calculate the number of the read processes of the data before update through the effect of the redundant data generating method selecting function 37. This number of the read processes is required to generate the redundant data through the effect of the method of read and modify. This is calculated as the number of the slots of the write data plus one. Next, the processor 30 operates to calculate the number of read processes through the effect of the redundant data generating method selecting function 37. This number of read processes is required to generate the redundant data through the effect of the method of all stripes. This is a value calculated by subtracting the number of the disk drives 300 for saving the write data from the number of the disk drives 300 for saving the data. Then, the redundant data generating method selecting function 37 is executed to compare the necessary number of reads in the method of all stripes with the necessary number of reads in the method of read and modify. If the conditional expression of:

$$\begin{aligned} & \text{(the necessary number of reads in the method of all stripes)} \leq \\ & \text{(the necessary number of reads in the method of read and modify)} \end{aligned} \quad (1)$$

is met, the method of all stripes is selected. Then, the operation goes to a step 1160.

The method of all stripes (second redundant data generating means) is made to be a stream of data shown in Fig. 6. At first, the request for read is issued to the disk drive 300 required for a read process, that is, the drive controller 50 for controlling the disk drives 300 (300d to 300f) for saving data that are not intended for the write process. Then, the data on the disk drives 300 (300d to 300f) are read by the read plane 112 of the cache memory 110 (see (1) of Fig. 6). The write data saved on the write plane 111 and the data read on the read plane 112 are transferred to the redundant data generator 130 (see (2) of Fig. 6). The redundant data generator 130 operates to generate the redundant data and saves it in the read plane 112 of the cache memory 100 (see (3) of Fig. 6). Next, the drive controller control function 36 utilizes the drive controller 60 so that the redundant data on the read plane 112 of the cache memory 110 is transferred to the disk drive 300 (300g) for saving the redundant data. This is the completion of reflecting the redundant data on the disk drive 300 (see (4) of Fig. 6).

In place, if the expression 1 is not met, the operation goes to a step 1030. As shown in (1) of Fig. 7 or 8, the request for reading the write data before update is issued to the disk drives 300 (300a to 300b) for saving the write data and then is transferred onto the read plane 112 of the cache memory 110. Then, the operation goes to a step 1040. The drive state reference function 33 is executed to check if the disk drive 300 (300g) for saving the redundant data is in use. If it is in use, the operation goes to a step 1170.

If the disk drive 300 (300g) for saving the redundant data is not in use, the operation goes to a step 1090. The redundant data generating method selecting function 37 is executed to derive a cylinder number (#) on the disk medium 350 from the position on the disk medium 350 where the redundant data is saved and calculate a data length at each turn of the cylinder # (track). In the case of a constant density recording whose information recording capacity per cylinder of an inner circumference is different from that of an outer circumference, the data length for one circumference can be easily obtained since the data length is a function of the cylinder #. In addition, if the information recording capacity per cylinder of the inner circumference is equal to that of the outer circumference, this calculation is not necessary. The data length per circumference in any cylinder can be immediately obtained from the device specification of the disk drive 300.

Next, the microprocessor 30 operates to compare the data length per circumference with the data length of the write data saved in the cache directory 120 and determine if the spinning on standby takes place. In actual, an overhead takes place for operating the redundant data inside of the micro program. Assuming that the overhead is $1/n$ of one spin time of the disk medium 350. Since the write data length is required to be within one spin time containing the overhead time, the condition for comparison is:

$$\text{Write Data Length} < (\text{Data Length per circumference of Disk Medium}) \times (1 - 1/n) \quad (\text{expression 2})$$

When this condition is met, the redundant data generating method selecting function 37 is executed to select the method for generating the redundant data through the effect of the method of a generation in a drive (third means for generating redundant data). Then, the operation goes to a step 1100. At a step 1090 shown in Fig. 5, it is assumed that an overhead ($1/n$) for operating the redundant data inside of the micro program is $1/4$.

The operation in this assumption is executed as shown in Fig. 7. The write data on the write plane 111 and the data

EP 0 844 561 A2

before update on the read plane 112 are transferred to the disk drive 300 (300g) (see (2) of Fig. 7). The write data transferred from the disk control device 20 and the data before update are transferred to a drive buffer 330 through a disk I/F 310. In synchronous with it, the microprocessor 320 located inside of the disk drive 300 operates to position the head at the position where the redundant data before update is saved through the disk control mechanism 340. On the termination of the positioning operation, the redundant data before update is read into the drive buffer 330 (see (3) of Fig. 7). The exclusive OR of the redundant data before update to be read, the write data having been stored in the buffer, and the data before update is calculated by the microprocessor 320 for generating the new redundant data and storing it in the drive buffer 330 (see (4) of Fig. 7). Next, when the disk medium 350 is spinned once and reaches the positioning point, the microprocessor 320 operates to write the redundant data on the drive buffer 330 (see (5) of Fig. 7). The disk drive 300 reports a success or failure of the generation and the write of the redundant data to the disk control device 20. When the write of the redundant data becomes successful, the operation goes to the next step 1110. The microprocessor 30 operates to detect if the update of the redundant data becomes successful or failed. If it is successful, the operation goes to a step 1120.

When the update of the redundant data is failed, the operation goes to a step 1170. At this step, for the retry process to be executed if the update of the redundant data is failed or the redundant data is in use, the redundant data generating method selecting function 37 is executed to select the method of read and modify.

The method of read and modify (first means for generating the redundant data) is executed as shown in Fig. 8. That is, since the data before update has been read into the read plane 112 of the cache memory 110, only the redundant data before update on the disk drive 300 (300g) for saving the redundant data is read onto the read plane 112 of the cache memory 110 (see (2) of Fig. 8). Then, the write data, the data before update, and the redundant data before update are transferred into the redundant data generator 130 (see (3) of Fig. 8) for generating the redundant data from those pieces of data and then saving the generated redundant data on the write plane 111 of the cache memory 110 (see (4) of Fig. 8). Next, the generated redundant data is written in the disk drive 300 (300g) for saving the redundant data (see (5) of Fig. 8). This is the completion of updating the redundant data.

Next, if neither of the conditional expressions (1) and (2) are met, the operation goes to a step 1140. At this step, the redundant data generating method selecting function 37 is executed to select the method of difference (fourth means for generating the redundant data) as the method for generating the redundant data. As shown in Fig. 9, the flow of the process is as shown in Fig. 9. The difference data generator control function 35 is executed to transfer the write data on the write plane 111, the data before on the read plane 112, and the difference data generator 140 (see (2) of Fig. 9) for generating the difference data and saving it on the read plane 112 (see (3) of Fig. 9) of the cache memory 110.

Assuming that the redundant data is the exclusive ORed data, for example, the difference data, termed herein, means the data given by taking an exclusive OR of the data A and B on the write plane 111 and the previous data a and b on the corresponding read plane 112 to the write plane 111, all of which are shown in Fig. 9.

Next, the drive controller control function 36 is executed to transfer the difference data on the read plane 112 of the cache memory 110 to the disk drive 300 (300g) for saving the redundant data through the use of the drive controller 50 (see (4) of Fig. 9). The difference data transferred to the disk control device 20 is then transferred to the drive buffer 330 through the disk I/F 310. In synchronous with it, the microprocessor 320 located inside of the disk drive 300 operates to position the head at the position where the redundant data before update is saved through the disk control mechanism 340. When the positioning operation is terminated, the microprocessor 320 operates to read the redundant data before update (see (5) of Fig. 9) and calculate an exclusive OR of the redundant data before update to be read and the difference data saved in the drive buffer 330 for generating the redundant data (see (6) of Fig. 9) and saving it in the drive buffer 330. Next, when the disk medium 350 spins once and reaches the positioning point, the microprocessor 320 operates to write the redundant data saved in the drive buffer 330 on the disk medium 350. If the write of the redundant data becomes successful, the disk drive reports the normal termination to the disk control device 20.

On the other hand, the microprocessor 30 located inside of the disk control device 20 operates to recognize that the generation of the redundant data becomes successful through the drive controller. Next, the operation goes to a step 1120. At this step, the write data is written in the disk drive 300. Then, proceeding to a step 1130, a series of writing operations are terminated.

The disk array device and the method for controlling the disk array device according to this embodiment are controlled to select the most approximate method for generating the redundant data from the method of all stripes and the method of read and modify, both of which are executed to generate the redundant data on the disk control device 20 according to the length of the write data received from the host 10, the access pattern such as sequential accesses, whether or not the redundant data is in use (that is, the magnitude of load) in the disk drive 300, and whether or not a failure takes place in the process of generating and saving the redundant data, as well as the method of a generation in a drive and the method of difference, both of which are executed to generate the redundant data on the disk drive 300. Hence, the disk array device and the control method thereof make it possible to reduce an overhead accompanied with the generation of the redundant data in processing the write data, improve a response time and a throughput of the

EP 0 844 561 A2

disk array device in the process for writing the data, and enhance the reliability of generation of the redundant data.

The foregoing description has been concretely expanded along the embodiments of the invention. However, it goes without saying that the present invention is not limited to the foregoing embodiments and may be modified into various modes without departing from the spirit of the invention.

5 For example, the present invention may widely apply to a magnetic disk device provided with the magnetic disk as a medium as well as a disk array device providing as a storage unit a disk drive having a general rotary storage medium such as an optical disk unit or a magnetooptical disk unit.

The disk array device according to the invention offers the effect of avoiding the spinning on standby caused if the data length of the write data sent from the host system is longer than one spin of the disk medium and thereby improving a response time when the redundant data is generated by the disk drive composing a disk array.

10 The disk array device according to the invention offers the effect of selecting the most approximate redundant data generating method that makes the necessary number of reads of the data before update minimal according to the divisional number of the write data received from the host system and thereby improving a throughput of the disk array device.

15 The disk array device according to the invention offers the effect of reducing an occupation time for one spin of the disk drive when the redundant data is generated by the disk drive composing the disk array and thereby improving the response time of the disk array device.

The disk array device according to the invention offers the effect of enhancing the efficiency of generating the redundant data accompanied with the processing of the write data according to the access pattern requested from the upper system and thereby improving a throughput of the disk array device.

20 The disk array device according to the invention offers the effect of enhancing the reliability of generating the redundant data in the case of generating the redundant data in the disk drive composing the disk array.

The method for controlling the disk array device according to the invention offers the effect of avoiding the spinning on standby caused if the data length of the write data sent from the host system is longer than one spin of the disk medium and thereby improving a response time on which the redundant data is generated by the disk drive composing the disk array.

25 The method for controlling the disk array device according to the present invention offers the effect of selecting the most approximate redundant data generating method in which the necessary number of reads of the data before update is made minimal according to the divisional number of the write data received from the host system and thereby improving a throughput of the disk array device.

30 The method for controlling the disk array device according to the present invention offers the effect of reducing an occupation time for one spin of the disk drive when the redundant data is generated by the disk drive composing the disk array and thereby improving the response time of the disk array device.

35 The method for controlling the disk array device according to the invention offers the effect of enhancing a generating efficiency of the redundant data accompanied with the processing of the write data and thereby improving a throughput of the disk array device.

The method for controlling the disk array device according to the invention offers the effect of enhancing the reliability of generating the redundant data when the redundant data is generated by the disk drive composing the disk array.

40 **Claims**

1. A disk array device comprising:

45 a plurality of disk drives (300) for saving data to be sent to a host system;
a subset of said disk drives composing a specific data group for saving write data received from said host system and at least one piece of redundant data generated from the write data;
a disk control device (20) for controlling data transfer between said host system and said disk drive (300);
a plurality of circuits (130) for generating said redundant data in the different methods from each other; and
50 a selective control circuit (30) for selecting at least one of said redundant data generating circuit (130) and executing the selected one.

2. The disk array device as claimed in claim 1, wherein said disk control device (20) operates to select a proper one of said redundant data generating circuits (130) based on at least one of a length of said write data received from said host system, an access mode specified by said host system, and a using status of said disk drive (300).

55 3. The disk array device as claimed in claim 2, wherein the selective control circuit (30) of said disk control device (20) selects such a redundant data generating circuit (130) as reducing a processing time of said write data containing generation of said redundant data to a minimum.

EP 0 844 561 A2

4. A disk array device comprising:

a plurality of disk drives (300) for saving data to be sent to a host system;
a subset of said disk drives (300) composing a specific data group for saving write data received from said host system and at least one piece of redundant data generated from said write data, each of said disk drives (300) containing a redundant data generating circuit for generating the redundant data from the data received from said host system and the data read from said disk;
a disk control device (20) for controlling data transfer between said host system and said disk drive (300);
said disk control device (20) having a redundant data generating circuit (130) and a selective control circuit (30) for selecting said redundant data generating circuit inside of said drive (300) or said redundant data generating circuit inside of said disk control device (20); and
said selective control circuit (30) serving to select the use of said redundant data generating circuit inside of said drive (300) only, the use of said redundant data generating circuits inside of said drive (300) and said control device (20), or the use of said redundant data generating circuit inside of said control device (20) only.

5. The disk array device as claimed in claim 4, wherein said disk control device (20) selects said first redundant data generating circuit or said second redundant data generating circuit based on a length of said write data received from said host system.

6. The disk array device as claimed in claim 4, wherein said disk control device (20) includes a plurality of said second redundant data generating circuits for generating said redundant data in different methods from each other.

7. A disk array device comprising:

a plurality of disk drives (300) for saving data to be sent to a host system;
a subset of said disk drives (300) composing a specific data group for saving divided parts of write data received from said host system and at least one piece of redundant data generated from the divided part of said write data in a distributed manner;
each of said disk drives (300) having a redundant data generating circuit for generating the redundant data from the data received thereby and the data read out of a disk;
a disk control device (20) for controlling data transfer between said host system and said disk drive (300),
said disk control device (20) including;
a first redundant data generating circuit inside of said control device (20) itself for generating new redundant data from said divided write data received from said host system, previous data to be updated by said divided write data, and redundant data of said divided write data about said data group;
a second redundant data generating circuit inside of said control device (20) for generating new redundant data from said divided write data received from said host system and data disabled to be updated by said divided write data about said data group;
a difference data generating circuit (140) for generating difference data from said divided write data received from said host system and the previous data to be updated by said divided write data; and
a selective control circuit (30) for selecting a circuit for generating said redundant data; and
wherein said selective control circuit operates to select for generating said redundant data the use of said first redundant data generating circuit inside of said control device (20), the use of said second redundant data generating circuit, the method of generating difference data generated by said difference data generating circuit (140), transferring said difference data to said disk drive (300) for saving the redundant data of said divided write data about said data group, and generating new redundant data from said difference data and redundant data with said redundant data generating circuit inside of said drive, or the method of transferring said divided write data received from said host system and the previous data to be updated by said divided write data to said disk drive (300) for saving the redundant data of said divided write data about said data group, and generating new redundant data from said write data, said previous data and said redundant data with said redundant data generating circuit inside of said drive (300).

8. The disk array device as claimed in claim 7, wherein

said selective control circuit (30) selects:
said second second redundant data generating means inside of said control device (20) if the conditional expression of

EP 0 844 561 A2

$b + 1 \leq 2 \cdot a + 1$ (a is a divisional number of said write data, b is a number of said disk drives (300) used except for saving the redundant data included in said data group)

5 is met;
said first redundant data generating means inside of said control device (20) if the drive for saving said redundant data is in use;
transferring said divided write data and the previous data to be updated by said divided write data to said disk drive for saving the redundant data of said write data about said data group and generating new redundant data from said write data and said previous data and said redundant data with said redundant data generating circuit inside of said drive (300) if said write data length is equal to or shorter than a given length; and
10 generating difference data generated by said difference data generating circuit, transferring said difference data to said disk drive for saving the redundant data of said divided write data about said data group, and generating new redundant data from said difference data and redundant data with said redundant data generating circuit inside of said drive (300).
15

9. The disk array device as claimed in claim 7, wherein the selective control circuit of said disk control device (20) selects such a redundant data generating circuit as reducing the processing time of said write data containing generation of the redundant data to a minimum.

20

25

30

35

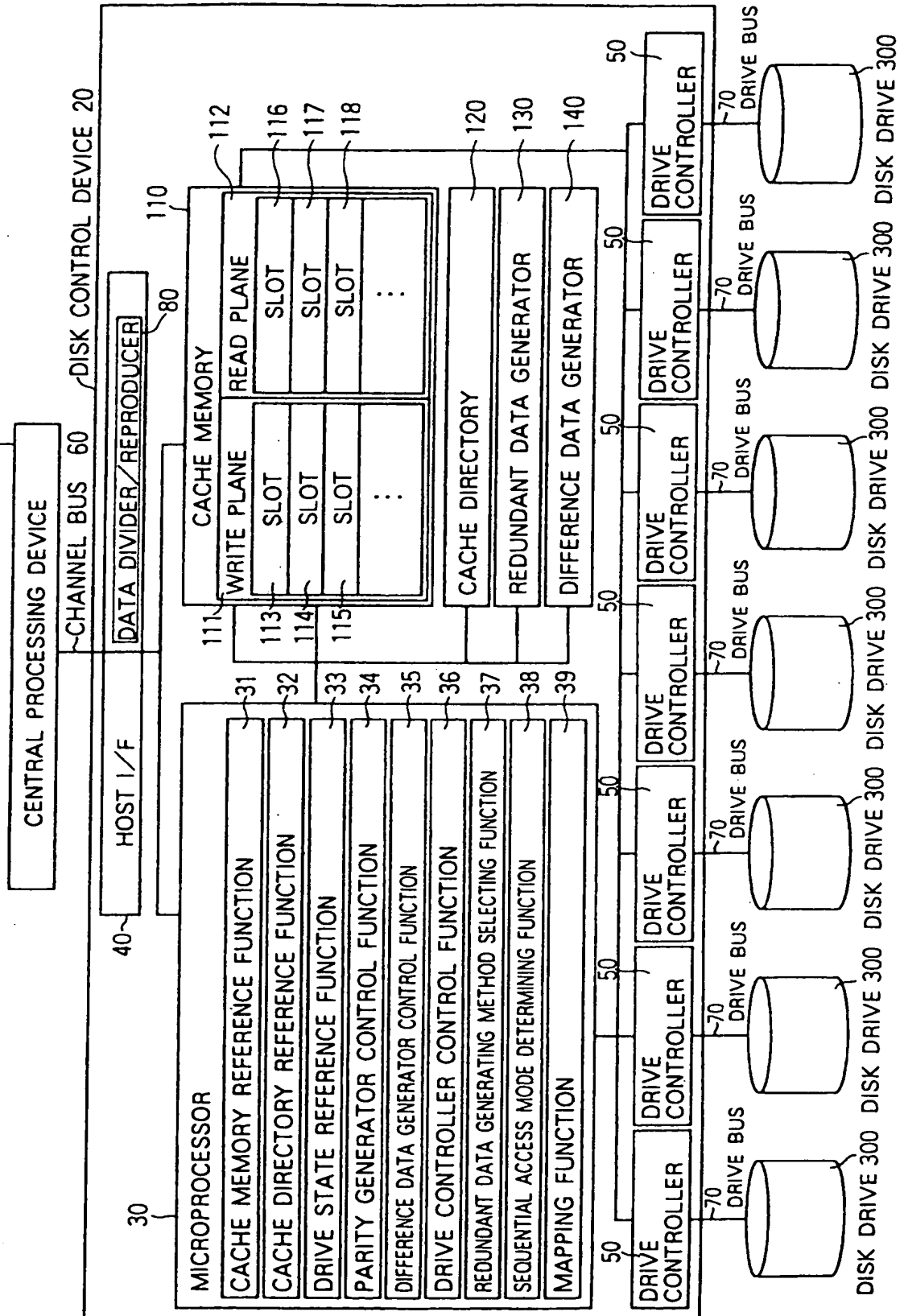
40

45

50

55

FIG. 1



EP 0 844 561 A2

FIG. 2

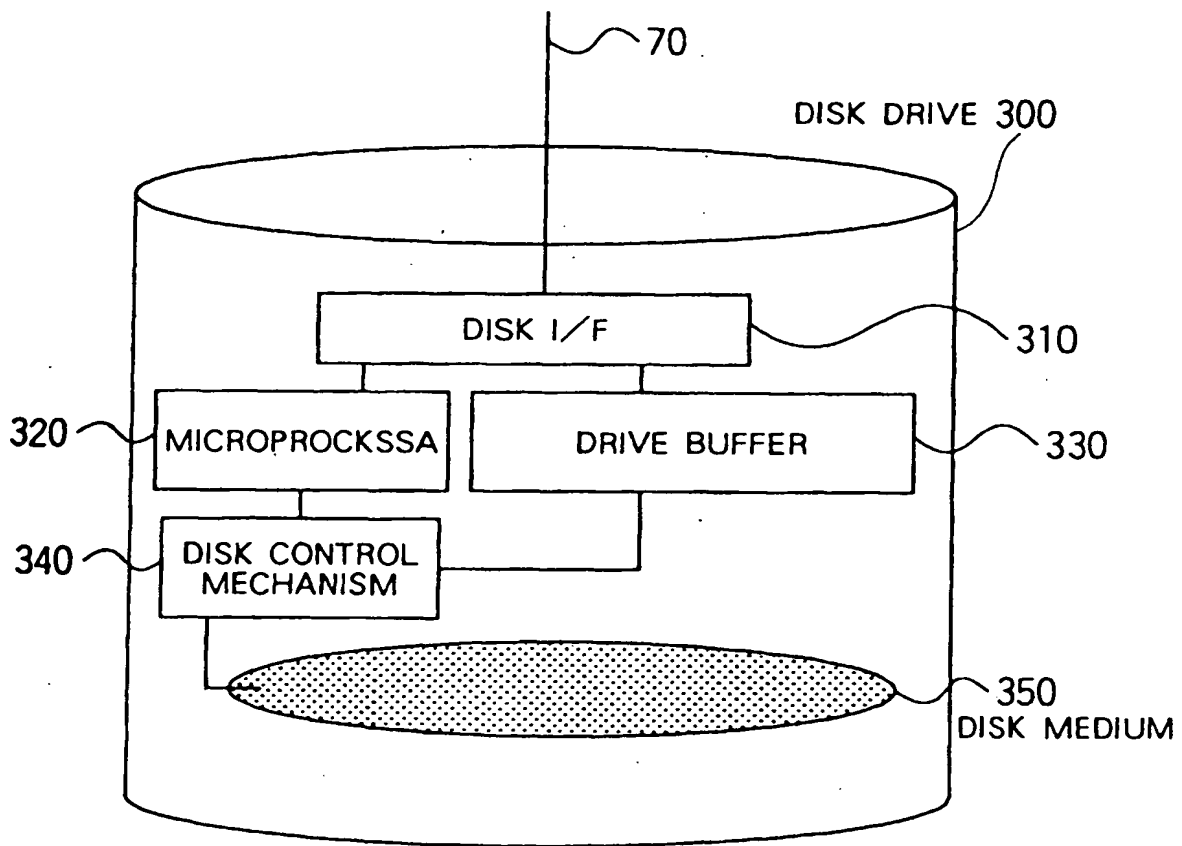


FIG. 3

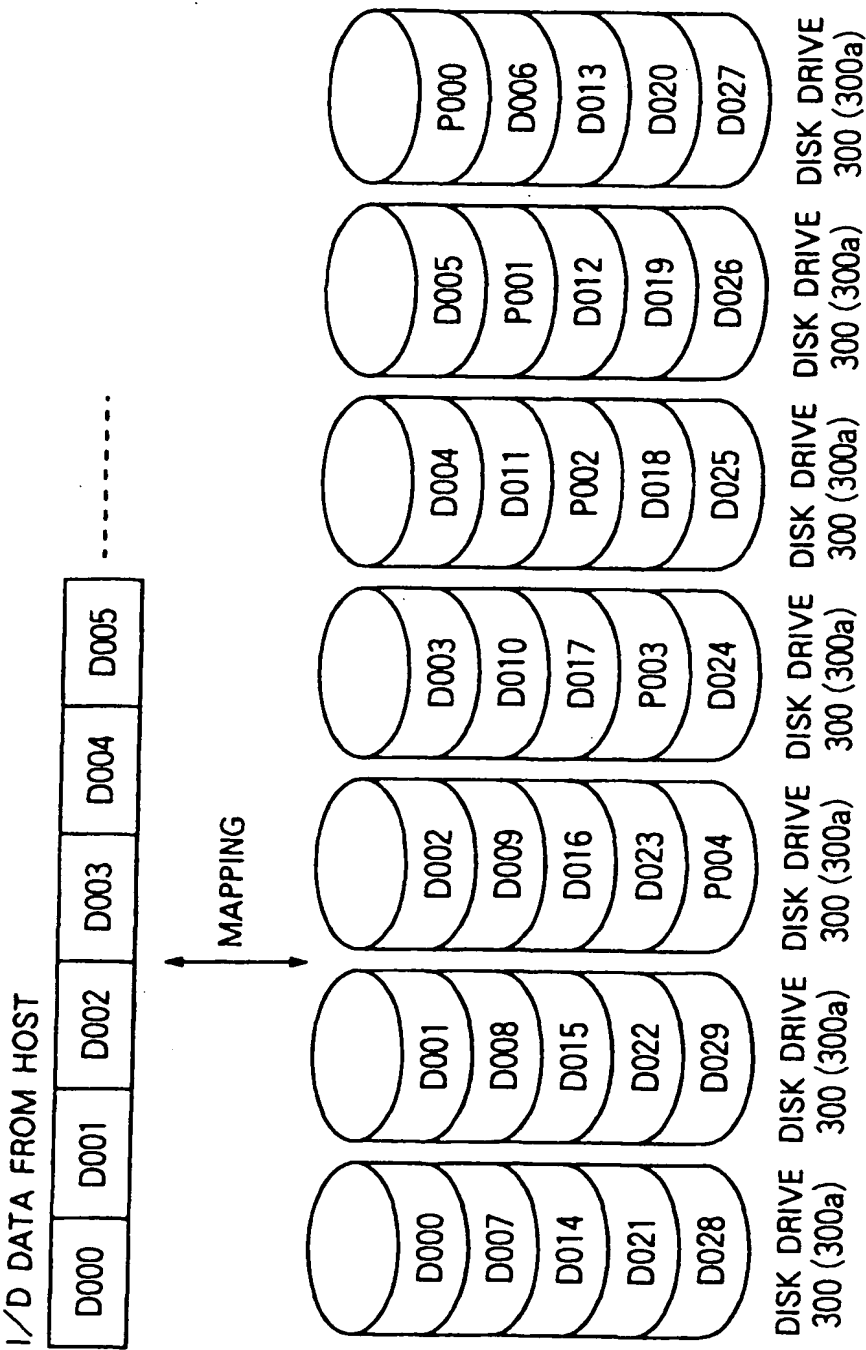
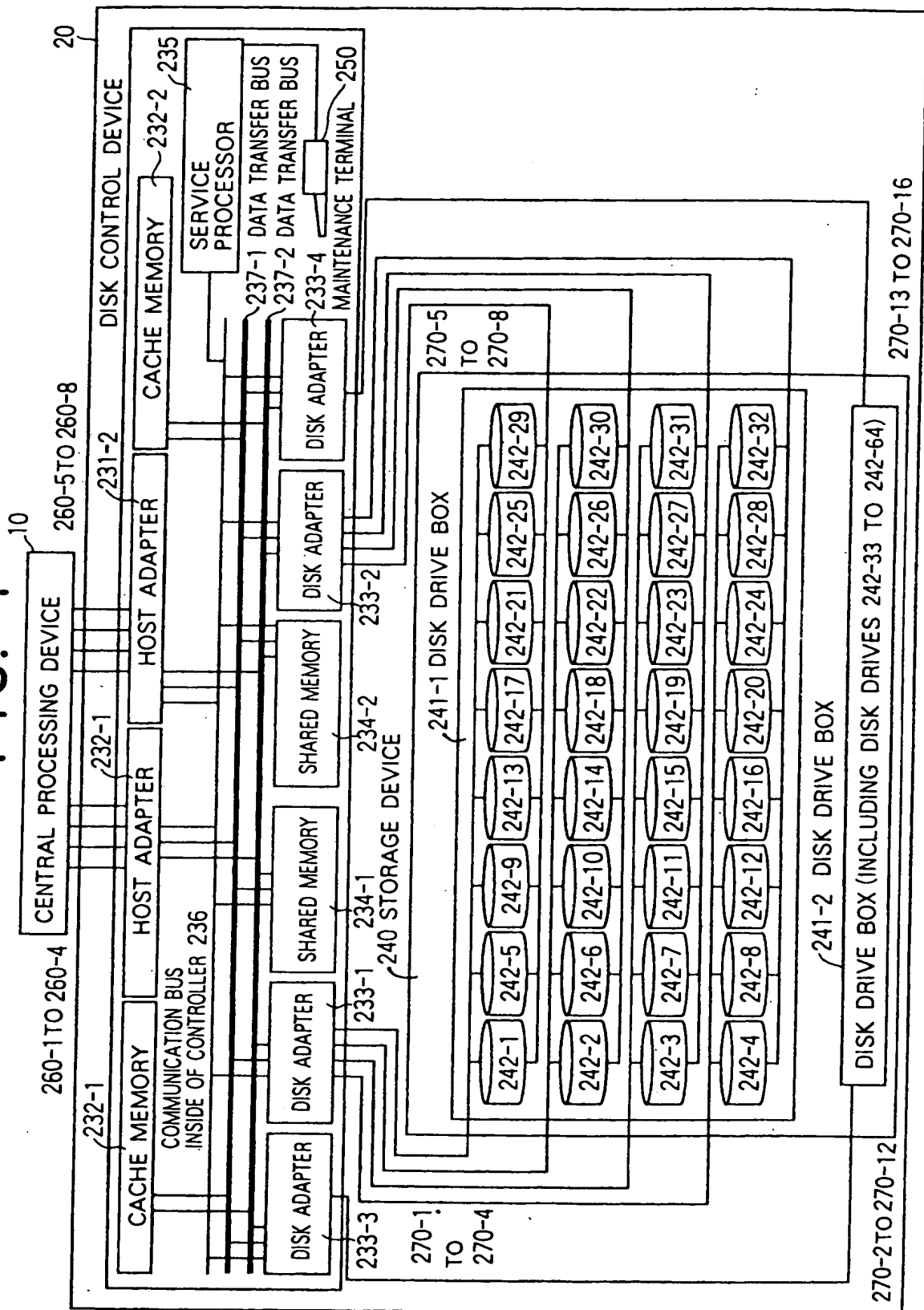
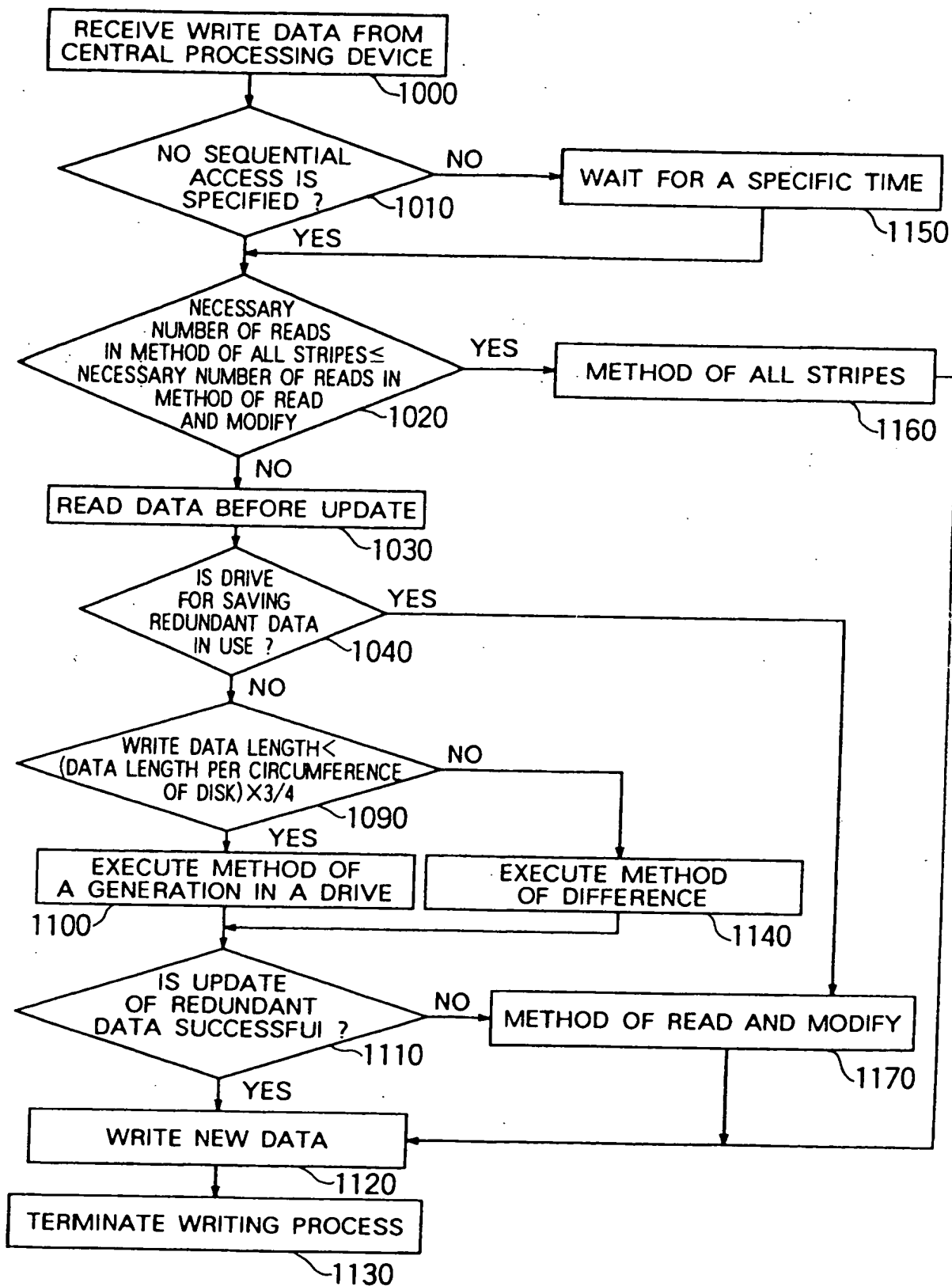


FIG. 4



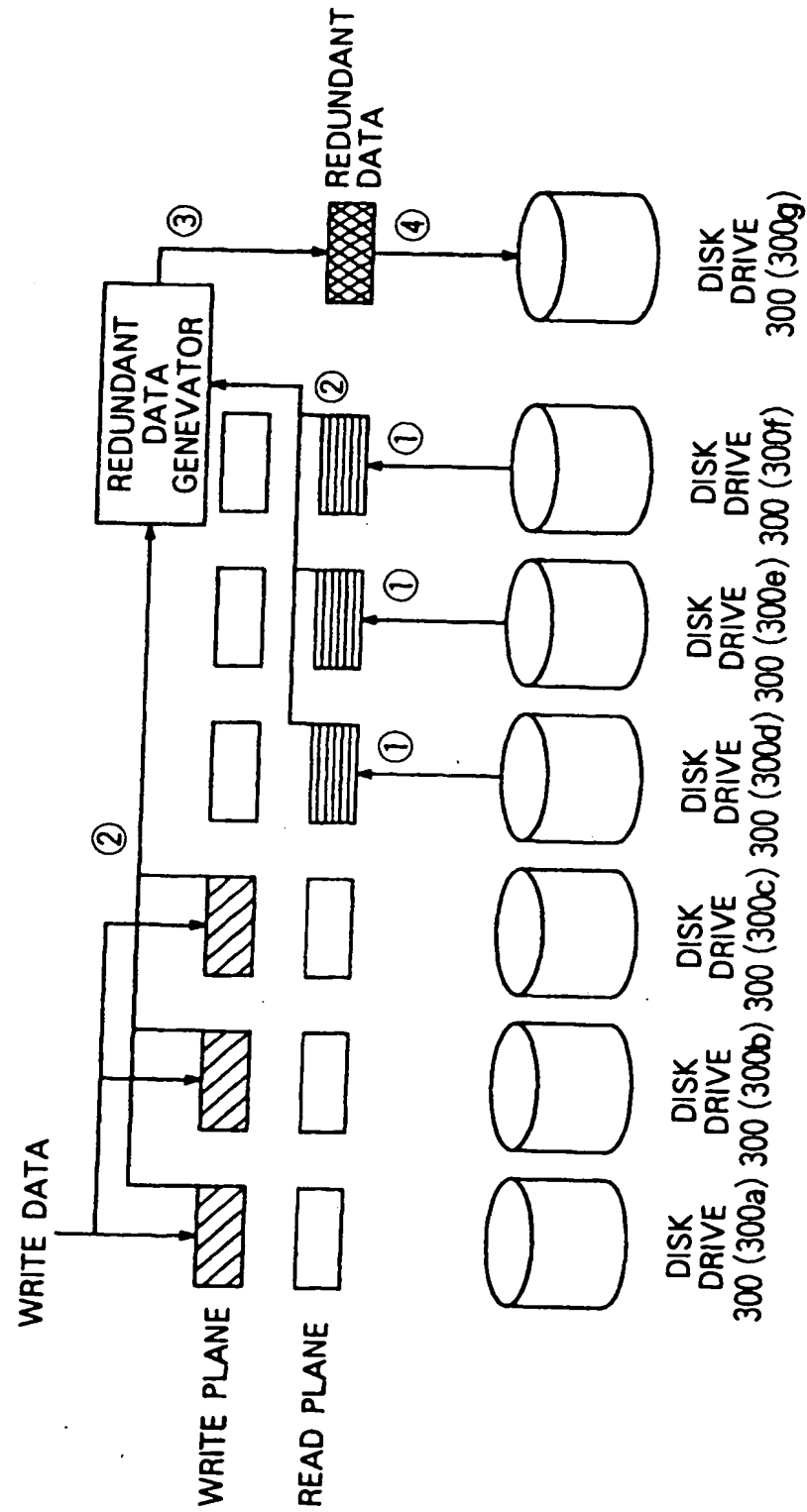
EP 0 844 561 A2

FIG. 5



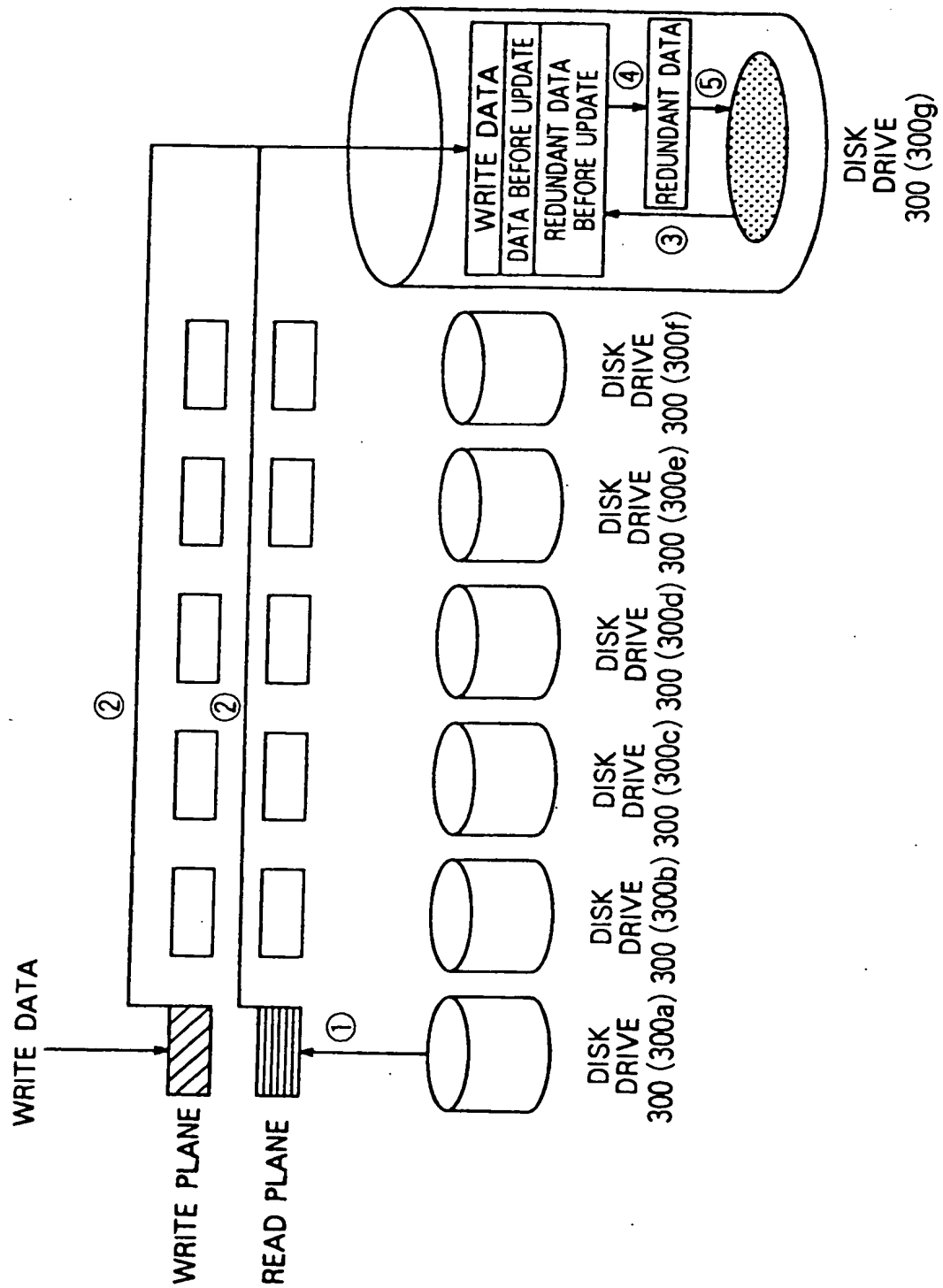
EP 0 844 561 A2

FIG. 6



EP 0 844 561 A2

FIG. 7



EP 0 844 561 A2

FIG. 8

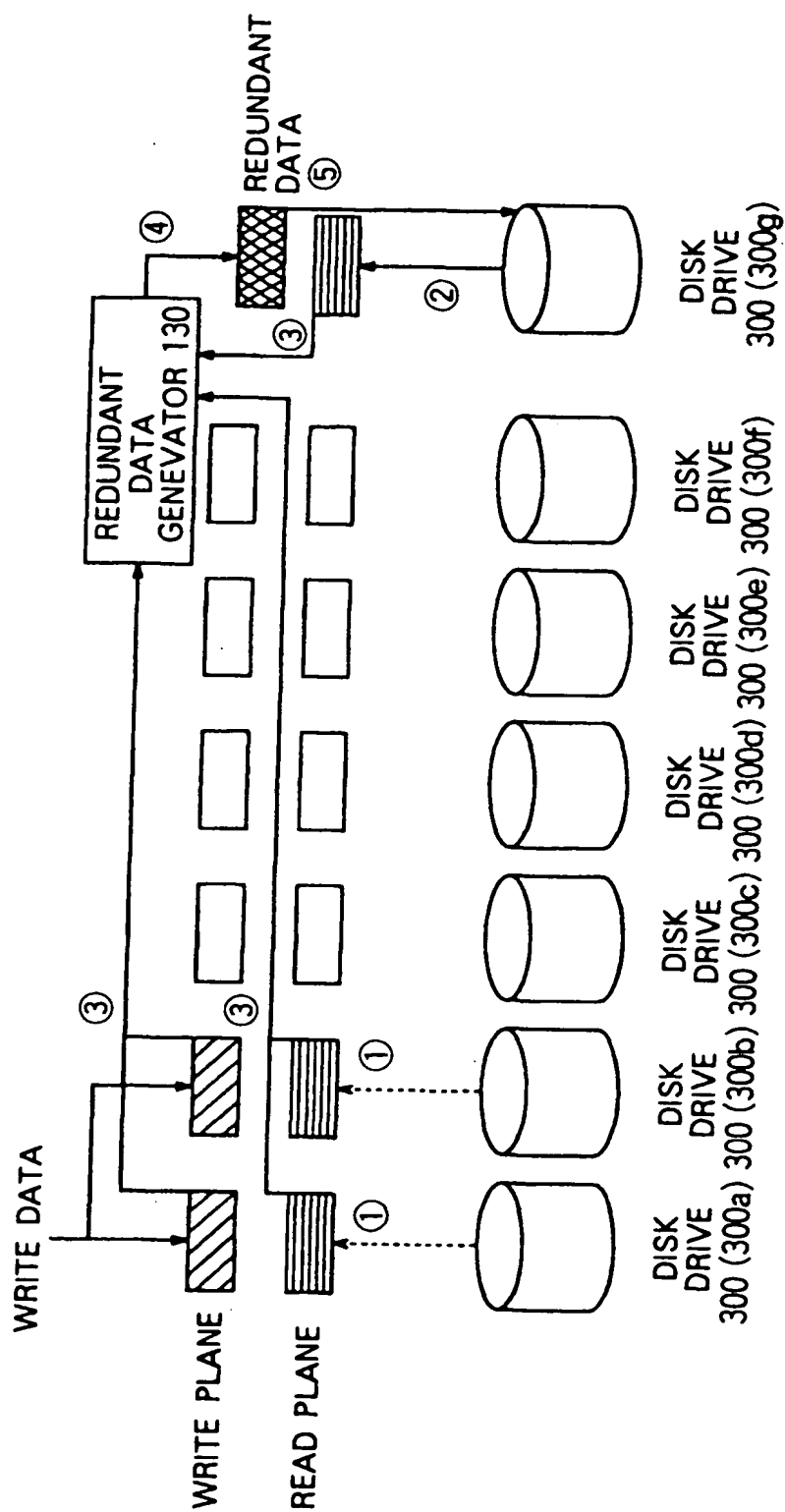


FIG. 9

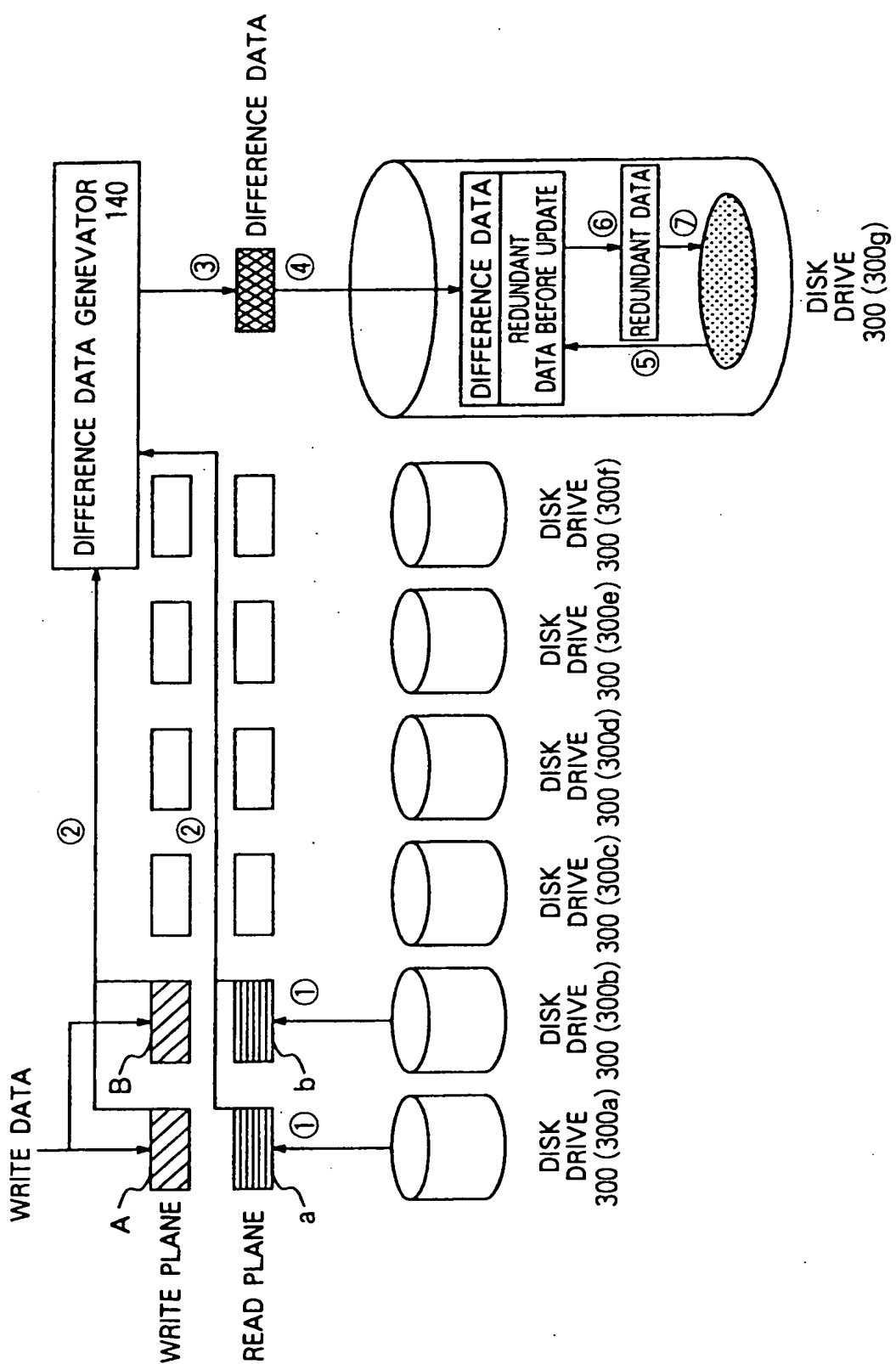


FIG. 10

